

REMARKS

The Office Action of January 24, 2003 has been carefully considered. Reconsideration of this application, as amended, is respectfully requested. Claims 1, 7-39, 41-43 and 45-49 are pending in this application. Of these, claims 1, 39, 43 and 49 are independent. In this Amendment, claims 1, 7-8, 10-19, 21-23, 27-29, 34, 39, 41-43, and 45-49 have been amended, inserting the word "digital" in front of "document" and "documents". Support in the specification for these amendments includes page 2, lines 11-29; page 10, lines 19-24; and page 15.

35 USC § 103

Claims 1, 7-22, 28, 39, 41-43, 45, 47 and 49 stand rejected under 35 U.S.C. § 103(a) as being unpatentable over Tanaka et al. ("Tanaka") in view of Pirolli et al. ("Pirolli"). The Office Action sets forth the rationale in support of these rejections in section 3, on pages 2-8. Applicants' respectfully request that these rejections be withdrawn for at least two reasons. First, Tanaka does not teach what it is alleged to teach. In particular, in contrast to all of the amended independent claims, Tanaka does not teach a method for processing digital documents. Tanaka processes hard copy documents, which he calls "image originals" at column 10, lines 49-59, in an effort to produce a digital document whose textual content closely resembles that of the image original. This process is known as Optical Character Recognition (OCR). This explains why Tanaka's patent is titled: Method and Apparatus for *Character Recognition* and why characters are illustrated in Figures 13A, 13B, 13C, 15, 19, 20A, 20B, 20C, 20D, 20E, 22A, 22B, 22C, 22D, 22E, 24A, 24B, 24C, 24D, 27, 28A, 28B, 31, 32, 33, 35, 37, 38, and 39. Once a document has been converted from a hard copy, or an image original, into a digital document Tanaka is finished with it. In contrast, applicants' claimed method is not applicable until a document is in a digital format. While Tanaka's method uses feature vectors, these differ from applicants' claimed feature vectors. Tanaka's feature vectors represent information about the characters that are being

C

recognized, while applicant's claimed feature vectors represent characteristics of digital documents. Second, applicants respectfully request that the rejections of claims 1, 7-22, 28, 39, 41-43, 45, 47 and 49 be withdrawn because there is no motivation to combine Tanaka and Pirolli. To establish a prima facie case of obviousness, there must be some suggestion or motivation, either in the references themselves or in the knowledge generally available to one of ordinary skill in the art, to modify the reference or to combine reference teachings. Applicants respectfully submit that Tanaka is not properly combined with Pirolli. As discussed above, Tanaka teaches recognizing characters on hard-copy documents, whereas Pirolli teaches analysis of linked digital documents. For these reasons, applicants respectfully request that the rejections of claims 1, 7-22, 28, 39, 41-43, 45, 47 and 49 be withdrawn and that amended claims 1, 7-22, 28, 39, 41-43, 45, 47 and 49 be allowed.

Allowable Subject Matter

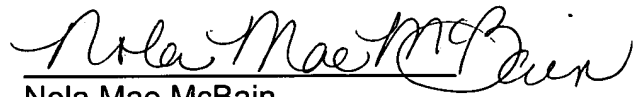
Claims 23-27, 29-38, 46 and 48 are objected to as being dependent upon a rejected based claim, but would be allowable if rewritten in independent form including all of the limitations of the base claim and any intervening claims. Applicants respectfully submit that claims 23-27, 29-38, 46 and 48 are in allowable condition because they depend upon allowable independent claims, as discussed herein above with respect to independent claims 1, 39, 43 and 49. Applicants therefore respectfully request that the objections to claims 23-27, 29-38, 46 and 48 be withdrawn and claims 23-27, 29-38, 46 and 48 be allowed.

Reconsideration/Admittance Requested

In view of the foregoing remarks and amendments, reconsideration of this application and allowance thereof are earnestly solicited.

In the event the Examiner considers personal contact advantageous to the disposition of this case, he is hereby authorized to call Applicant's attorney, Nola Mae McBain, at Telephone Number (650) 812-4264, Palo Alto, California.

Respectfully submitted,

A handwritten signature in cursive script, reading "Nola Mae McBain". The signature is written in dark ink and is positioned above a horizontal line.

Nola Mae McBain
Attorney for Applicant(s)
Registration No. 35,782
Telephone: 650-812-4264

Date: April 24, 2003

C

APPENDIX A

Marked Up Amended Claims Under 37 C.F.R. 1.121(c)(1)(ii):

Appendix A sets forth a marked up version of the prior pending amended claims for their corresponding pending claims with additions shown with underlining (e.g. new text) and deletions shown with a strikethrough (e.g. ~~delete text~~).

1. (Thrice Amended) A method for quantitatively representing digital documents in a vector space, comprising the steps of:

identifying a first digital document to be processed from a plurality of digital documents;

extracting a first feature corresponding to the first document from the plurality of digital documents, the first feature comprising text surrounding an image included in the digital document, the text surrounding the image not being anchor text;

converting the first feature to a first vector; and

associating the first vector with the first digital document.

7. (Amended) The method of claim 1 further comprising the steps of:

extracting a second feature corresponding to the digital document, the second feature comprising a first URL representing the first digital document;

converting the second feature to a second vector; and

associating the second vector with the first digital document.

8. (Amended) The method of claim 7, wherein the step of converting the second feature comprises the sub-steps of:

identifying each unique word within the URLs representing all digital documents in the collection of digital documents; and

counting the occurrences of each unique word in the first URL;

creating a vector having a number of dimensions equal to the number of unique words in the URLs representing all digital documents in the collection of digital documents, and further having as each element a numeric value representative of the number of occurrences in the first URL of the corresponding word.

10. (Amended) The method of claim 1 further comprising the steps of:
extracting a second feature corresponding to the first digital document, the second feature comprising inlinks in the collection of digital documents linking to the first document;

converting the second feature to a second vector; and
associating the second vector with the first digital document .

11. (Amended) The method of claim 10, wherein the step of converting the second feature comprises the sub-steps of:

identifying each digital document having links within the collection of digital documents;

determining how many times each digital document having links points to the first digital document; and

creating the second vector having a number of dimensions equal to the number of digital documents having links in the collection of digital documents, and the second vector further having as each element a numeric value representative of the number of links in each corresponding digital document linking to the first digital document.

12. (Amended) The method of claim 11, wherein the numeric value representative of the number of links in each corresponding digital document linking to the first digital document is calculated as the token frequency weight of the corresponding link multiplied by the inverse context frequency weight of the corresponding link.

13. (Amended) The method of claim 10, wherein the step of converting the second feature comprises the sub-steps of:

identifying each digital document having hyperlinks within the collection of digital documents, and further identifying each unique word associated with URLs defining hyperlinks in each digital document;

counting the occurrences of each unique word in the URLs defining hyperlinks pointing to the first digital document; and

creating the second vector having a number of dimensions equal to the number of unique words associated with URLs defining hyperlinks within the collection of digital documents, and the second vector further having as each element a numeric value representative of the number of occurrences in the URLs defining hyperlinks pointing to the first digital document of the corresponding word.

14. (Amended) The method of claim 13, wherein the numeric value representative of the number of occurrences in the URLs defining hyperlinks pointing to the first digital document of the corresponding word is calculated as the token frequency weight of the corresponding word multiplied by the inverse context frequency weight of the corresponding word.

15. (Amended) The method of claim 1 further comprising the steps of:

extracting a second feature corresponding to the first digital document, the second feature comprising outlinks in the collection of digital documents linking to the first digital document;

converting the second feature to a second vector; and

associating the second vector with the first digital document.

16. (Amended) The method of claim 15, wherein the step of converting the second feature comprises the sub-steps of:

identifying each other digital document linked to by all digital documents within the collection of digital documents; and

C

creating the second vector having a number of dimensions equal to the number of other digital documents linked to by digital documents in the collection of digital documents, and the second vector further having as each element a numeric value representative of the number of links in the first digital document linking to each corresponding other digital document.

17. (Amended) The method of claim 16, wherein the numeric value representative of the number of links in the first digital document linking to each corresponding other digital document is calculated as the token frequency weight of the corresponding link multiplied by the inverse context frequency weight of the corresponding link.

18. (Amended) The method of claim 15, wherein the step of converting the second feature comprises the sub-steps of:

identifying each unique word associated with URLs defining hyperlinks in each digital document in the collection of digital documents;

counting the occurrences of each unique word in the URLs defining hyperlinks in the first digital document; and

creating the second vector having a number of dimensions equal to the number of unique words associated with the URLs defining hyperlinks in each digital document, and the second vector further having as each element a numeric value representative of the number of occurrences in the URLs defining hyperlinks in the first digital document of the corresponding word.

19. (Amended) The method of claim 18, wherein the numeric value representative of the number of occurrences in the URLs defining hyperlinks in the first digital document of the corresponding word is calculated as the token frequency weight of the corresponding word multiplied by the inverse context frequency weight of the corresponding word.

C

21. (Amended) The method of claim 20, wherein the step of converting the second feature comprises the sub-steps of:

for each possible text genre, processing the first digital document to calculate the probability that the first digital document is of the corresponding text genre; and

creating the second vector having a number of dimensions equal to the number of possible text genres, and the second vector further having as each element a numeric value representative of the probability that the first digital document is of the corresponding genre.

22. (Thrice Amended) The method of claim 49, wherein the first feature comprises the color histogram for the image included in the first digital document.

23. (Amended) The method of claim 22, wherein the step of converting the first feature comprises the sub-steps of:

quantizing the image represented by the first digital document into a multi-dimensional color model;

creating a color histogram having a plurality of bins for each dimension in the color model, each bin corresponding to a unique combination of binary bits representing information from the associated dimension of the color model;

counting each of a plurality of pixels from the image in a corresponding bin associated with each dimension of the color model; and

creating the first vector having a number of dimensions equal to the total number of bins in the color histogram, and the first vector further having as each element a numeric value representative of the number of pixels in the image corresponding to the corresponding histogram bin.

27. (Amended) The method of claim 23, wherein the image represented by the first digital document comprises a region of a bitmap.

C

28. (Amended) The method of claim 49, wherein the first feature comprises the color complexity of an image represented by the first digital document.

29. (Amended) The method of claim 28, wherein the step of converting the first feature comprises the sub-steps of:

- quantizing the image represented by the first digital document into a multi-dimensional color model;

- determining the maximum number of pixels in any row in any image represented by any digital document in the collection of digital documents;

- determining the maximum number of pixels in any column in any image represented by any digital document in the collection of digital documents;

- creating a horizontal complexity histogram and a vertical complexity histogram, each having a number of bins equal to the maximum number of pixels in any row and in any column, respectively;

- identifying horizontal runs of pixels of all possible lengths in the quantized image, and for each possible length, counting the number of pixels in a plurality of rows of the quantized image belonging to the horizontal runs in a corresponding bin of the horizontal complexity histogram;

- identifying vertical runs of pixels of all possible lengths in the quantized image, and for each possible length, counting the number of pixels in a plurality of columns of the quantized image belonging to the vertical runs in a corresponding bin of the horizontal complexity histogram;

- creating a horizontal complexity vector having a number of dimensions equal to the maximum number of pixels in any row, and further having as each element a numeric value representing the number of pixels in the image in the corresponding horizontal histogram bin; and

- creating a vertical complexity vector having a number of dimensions equal to the maximum number of pixels in any column, and further having as each element a numeric value representing the number of pixels in the image in the corresponding vertical histogram bin.

34. (Amended) The method of claim 28, wherein the step of converting the first feature comprises the sub-steps of:

quantizing the image represented by the first digital document into a multi-dimensional color model;

determining the maximum number of pixels in any row in any image represented by any digital document in the collection of digital documents;

determining the maximum number of pixels in any column in any image represented by any digital document in the collection of digital documents;

creating a horizontal complexity histogram and a vertical complexity histogram, each having a selected number of bins corresponding to a plurality of quantized ranges of run lengths;

identifying horizontal runs of pixels of all possible lengths in the quantized image, and for each possible length, counting the number of pixels in a plurality of rows of the quantized image belonging to the horizontal runs in a corresponding bin of the horizontal complexity histogram;

identifying vertical runs of pixels of all possible lengths in the quantized image, and for each possible length, counting the number of pixels in a plurality of columns of the quantized image belonging to the vertical runs in a corresponding bin of the horizontal complexity histogram;

creating a horizontal complexity vector having a number of dimensions equal to the selected number of bins in the horizontal complexity histogram, and further having as each element a numeric value representing the number of pixels in the image in the corresponding horizontal histogram bin; and

creating a vertical complexity vector having a number of dimensions equal to the number of bins in the vertical complexity histogram, and further having as each element a numeric value representing the number of pixels in the image in the corresponding vertical histogram bin.

39. (Thrice Amended) A signal representing instructions for quantitatively representing in a vector space users of a collection of digital documents, the instructions comprising:

identifying a first user to be processed from the users of the collection of digital documents;

extracting from the collection of digital documents a first feature representing a first sub-set of digital documents of the collection that have been accessed by the first user;

converting the first feature to a first vector; and

associating the first vector with the first user.

41. (Thrice Amended) The signal of claim 39, wherein the converting instruction comprises:

identifying each unique digital document in the collection of digital documents;

calculating the number of times the first user accessed each digital document in the collection of digital documents; and

creating the first vector having a number of dimensions equal to the number of digital documents in the collection of digital documents, and the first vector further having as each element a numeric value representative of the number of times the first user has accessed the corresponding digital document.

42. (Thrice Amended) The signal of claim 41, wherein the value representative of the number of times the first user has accessed the corresponding digital document is calculated as the token frequency weight of the corresponding digital document multiplied by the inverse context frequency weight of the corresponding digital document.

43. (Twice Amended) A computer-readable medium containing instructions for causing a computer-system to quantitatively ~~representing~~ represent digital documents in a vector space, by the steps of:

identifying a digital document to be processed from a plurality of digital documents;

selecting an image feature as a first feature, the image feature being associated with the non-text content of an image included in the digital document;

extracting from the document information associated with the first feature;

converting information associated with the first feature into a first vector;

associating the first vector with the digital document;

selecting a second feature from a set of multi-modal features including a user information feature and a genre feature;

extracting from the document information associated with the second feature;

converting the information associated with the second feature into a second vector; and

associating the second vector with the digital document.

45.(Twice Amended) The computer-readable medium of claim 43 wherein the first feature comprises a color histogram for the image included in the digital document.

46.(Amended) The computer-readable medium of claim 45 wherein converting the information associated with the first feature into the first vector comprises the steps of:

quantizing the image included in the digital document into a multi-dimensional color model;

creating a color histogram having a plurality of bins for each dimension in the color model, each bin corresponding to a unique combination of binary bits representing information from the associated dimension of the color model;

counting each of a plurality of pixels from the image in a corresponding bin associated with each dimension of the color model; and

C

creating a vector having a number of dimensions equal to the total number of bins in the color histogram, and further having as each element a numeric value representative of the number of pixels in the image corresponding to the corresponding histogram bin.

47. (Twice Amended) The computer-readable medium of claim 43 wherein the first feature comprises color complexity of the image included in the digital document.

48. (Amended) The computer-readable medium of claim 47 wherein converting the information associated with the first feature into the first vector comprises the steps of:

- quantizing the image included in the digital document into a multi-dimensional color model;

- determining the maximum number of pixels in any row in any image represented by any digital document in the collection of digital documents;

- determining the maximum number of pixels in any column in any image represented by any digital document in the collection of digital documents;

- creating a horizontal complexity histogram and a vertical complexity histogram, each having a number of bins equal to the maximum number of pixels in any row and in any column, respectively;

- identifying horizontal runs of pixels of all possible lengths in the quantized image, and for each possible length, counting the number of pixels in a plurality of rows of the quantized image belonging to the horizontal runs in a corresponding bin of the horizontal complexity histogram;

- identifying vertical runs of pixels of all possible lengths in the quantized image, and for each possible length, counting the number of pixels in a plurality of columns of the quantized image belonging to the vertical runs in a corresponding bin of the horizontal complexity histogram;

- creating a horizontal complexity vector having a number of dimensions equal to the maximum number of pixels in any row, and further having as each

C

element a numeric value representing the number of pixels in the image in the corresponding horizontal histogram bin; and

creating a vertical complexity vector having a number of dimensions equal to the maximum number of pixels in any column, and further having as each element a numeric value representing the number of pixels in the image in the corresponding vertical histogram bin.

49. (Twice Amended) A method for quantitatively representing digital documents in a vector space, comprising the steps of:

identifying a first digital document to be processed from a plurality of digital documents;

extracting a first feature corresponding to the first digital document from the plurality of digital documents, the first feature comprising an image feature associated with non-text content of an image included in the first digital document;

converting the first feature to a first vector;

associating the first vector with the first digital document;

extracting a second feature corresponding to the digital document, the second feature comprising a one of a user feature and a text genre feature;

converting the second feature into a second vector; and

associating the second vector with the first digital document.